



Paper Type: Original Article

Cultivating Clean Skies: Unveiling the Tapestry of Air Quality in Gujarat through Innovative Machine Learning Analysis

Gaddam Advitha¹, Allada Nagasai Varaprasad¹, Koti Vennela Khushi¹, Pullabhotla Vijay¹, Sukanta Nayak^{2,*} 

¹ School of Computer Science and Engineering, VIT–AP University, Inavolu, Beside AP Secretariat, Amaravati AP, India; advithagaddam30@gmail.com; nagasaiallada@gmail.com; vennelakhushi2004@gmail.com; vijaybharat0305@gmail.com.

² Department of Mathematics, School of Advanced Sciences, VIT – AP University, Inavolu, Beside AP Secretariat, Amaravati AP, India; sukantgacr@gmail.com.

Citation:

Received: 20 August 2024
Revised: 23 October 2024
Accepted: 14 November 2024

Advitha, G., Khushi, K. V., Vijay, P., & Nayak, S. (2024). Cultivating clean skies: unveiling the tapestry of air quality in Gujarat through innovative machine learning analysis. *Big data and computing visions*, 4(4), 326-339.

Abstract


Air pollution emerges as a formidable threat to both public health and environmental integrity, especially in regions undergoing rapid development. In this context, Gujarat, situated in Western India, grapples with escalating air quality degradation attributable to industrialization, vehicular emissions, and agricultural practices. The imperative to accurately forecast Air Quality Indices (AQIs) becomes paramount for the prompt implementation of mitigation measures, thereby safeguarding public well-being. This research delves into the utilization of Machine Learning (ML) algorithms, specifically Random Forest (RF) and XGBoost, to predict AQIs in Gujarat. Four pivotal parameters, PM_{2.5}, PM₁₀, SO₂, and NO_x, are scrutinized due to their substantial impact on air quality and inclusion in publicly available datasets. Remarkably, both RF and XGBoost models exhibit outstanding performance, surpassing 99% accuracy. This exceptional capability underscores the transformative potential of ML in addressing the complex challenges posed by air pollution. Leveraging the precise predictions of AQI values, these models can catalyze the development of robust early warning systems and guide policymaking endeavors directed at enhancing air quality. Notably, this investigation unveils PM_{2.5} and PM₁₀ as primary culprits influencing AQI levels in Gujarat, underscoring the urgency for stringent emission control measures targeting these pollutants. Furthermore, the study sheds light on the significant impact of meteorological factors, especially maximum temperature, on AQI fluctuations, necessitating adaptive strategies to counteract climate change's repercussions on air quality.

Keywords: Air quality, Machine learning, Random forest, XGBoost, AQI prediction.

1 | Introduction

In the realm of environmental concerns, air pollution stands as a formidable challenge, posing severe threats to public health and ecosystem integrity. This complex mixture of particulate matter and gaseous pollutants,

 Corresponding Author: sukantgacr@gmail.com

 <https://doi.org/10.22105/bdcv.2024.485515.1213>



Licensee System Analytics. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0>).

largely stemming from human activities, has been linked to a range of respiratory ailments, cardiovascular complications, and adverse neurological impacts [1]. The escalating severity of air pollution has spurred a surge of research efforts, seeking to unravel the intricate mechanisms of this phenomenon and develop effective mitigation strategies [2].

Despite the extensive body of knowledge accumulated through past research, traditional statistical methods employed for air quality prediction have often exhibited limitations in capturing the intricate relationships between air pollutants and meteorological factors. This has fueled the exploration of alternative approaches, particularly Machine Learning (ML) techniques, which hold immense promise for extracting meaningful insights from large datasets and identifying subtle patterns that may elude conventional methods. Motivated by this potential, the present study delves into the application of two powerful ML algorithms, namely RF and XGBoost, for air quality prediction in Gujarat, India. This industrially thriving state in Western India has witnessed a surge in air pollution levels, necessitating accurate and timely predictions to inform mitigation measures and safeguard public health [3].

To address this pressing issue, we employ RF and XGBoost models trained on a comprehensive dataset [4] encompassing various parameters influencing air quality, including meteorological factors, pollutant concentrations, and temporal factors. The models' remarkable performance, achieving accuracies exceeding 99% in predicting Air Quality Indices (AQIs), underscores the efficacy of ML in addressing air pollution challenges. Our study extends beyond AQI prediction, providing valuable insights into the factors influencing air quality in Gujarat. We uncover the significant impact of PM_{2.5} and PM₁₀ on AQI levels, highlighting the need for stringent emission control measures targeting these pollutants.

Additionally, we emphasize the influence of meteorological factors, particularly maximum temperature, on AQI variations, necessitating adaptation strategies to counteract climate change's repercussions on air quality. Our research stands at the forefront of air quality research in Gujarat, offering a robust and effective ML-based approach to AQI prediction. Its findings have significant implications for policy decisions and mitigation strategies aimed at improving air quality and safeguarding public health in the region. Air pollution, a ubiquitous environmental hazard, refers to the presence of harmful substances in the atmosphere that degrade air quality and pose severe health risks. These substances originate from various sources, including natural processes and human activities, and can take the form of gaseous pollutants, particulate matter, and volatile organic compounds.

2| Factors Contributing to Air Pollution and Its Measuring Index

Air pollution, a ubiquitous environmental hazard, refers to the presence of harmful substances in the atmosphere that degrade air quality and pose severe health risks [5]. These substances originate from various sources, including natural processes and human activities, and they can take the form of gaseous pollutants, particulate matter, and volatile organic compounds.

- I. Industrial processes: combustion of fossil fuels in industries releases pollutants like SO₂, NO_x, and CO.
- II. Vehicular emissions: gasoline-powered vehicles emit pollutants like CO, NO_x, and hydrocarbons.
- III. Agricultural activities: burning of agricultural waste and use of fertilizers release pollutants like PM and ammonia.
- IV. Household activities: cooking and burning household waste release pollutants like PM and CO.

2.1| Air Quality Index

The AQI serves as a standardized measure of air quality, providing a simple and easily understandable representation of the overall air quality at a particular location and time. It is calculated based on the concentrations of five criteria pollutants viz. PM_{2.5}, PM₁₀, NO₂, SO₂, and CO [5]. AQI values are categorized into six color-coded levels, each representing a different level of health concern:

- I. Good (0-50): minimal or no health impact.
- II. Moderate (51-100): may cause mild health problems for sensitive individuals.
- III. Unhealthy for sensitive groups (101-150): may cause health problems for sensitive individuals and generally uncomfortable for everyone.
- IV. Unhealthy (151-200): may cause breathing problems for people with heart or lung disease and generally uncomfortable for everyone.
- V. Very unhealthy (201-300): may cause respiratory problems for everyone and may trigger serious health problems for people with heart or lung disease.
- VI. Hazardous (301-500): may cause serious health problems for everyone.

2.2 | Calculating AQI

The AQI is a numerical scale used to communicate how polluted the air currently is or how polluted it is forecast to become. The AQI is typically calculated based on concentrations of various air pollutants. The specific pollutants considered may vary, but common ones include particulate matter (PM10 and PM2.5), ground-level ozone (O3), sulphur dioxide (SO2), nitrogen dioxide (NO2), and carbon monoxide (CO).

The AQI is usually calculated based on the concentration of individual pollutants, and the final AQI value is determined by the highest sub-index among these pollutants. The general equation for calculating the sub-index for a specific pollutant is as follows:

$$I_i = \frac{(C_i - I_{low}) * (B_{high} - B_{low})}{(I_{high} - I_{low})} + B_{low},$$

where

- I. I_i is the sub-index for the individual pollutant.
- II. C_i is the concentration of the individual pollutant.
- III. I_{low} and I_{high} are the breakpoints that define the sub-index ranges for the pollutant.
- IV. B_{low} and B_{high} are the corresponding AQI values associated with the breakpoints.

Once the sub-indices for all relevant pollutants are calculated, the overall AQI is determined by selecting the maximum sub-index.

The equation for AQI is as follows:

$$AQI = \max(I_1, I_2, \dots, I_n),$$

where

- I. AQI is the overall AQI.
- II. I_1, I_2, \dots, I_n are the sub-indices for individual pollutants.

2.3 | Impact of Air Pollution in Gujarat and Larger Indian Context

Air pollution poses a significant threat to public health and the environment in India and Gujarat. The country's rapid industrialization, urbanization, and population growth have contributed to increasing air pollution levels [6].

In India, air pollution is estimated to cause over 2 million premature deaths annually [7]. Gujarat, a state in Western India, has experienced deteriorating air quality due to factors such as industrial emissions, vehicular traffic, and agricultural activities. Ahmedabad, the state's largest city, has consistently ranked among the most polluted cities in India [8].

Air pollution in India and Gujarat has far-reaching consequences, including:

- I. Respiratory illnesses: exposure to air pollutants can cause respiratory problems like asthma, bronchitis, and Chronic Obstructive Pulmonary Disease (COPD).
- II. Cardiovascular diseases: air pollution can exacerbate cardiovascular conditions like heart failure and arrhythmias [9].
- III. Premature mortality: air pollution is associated with an increased risk of premature death due to respiratory and cardiovascular complications.
- IV. Environmental degradation: air pollution contributes to acid rain, smog formation, and climate change [13].

Effective mitigation strategies are urgently needed to address the air pollution crisis in India and Gujarat. These strategies should focus on reducing emissions from industries, vehicles, and agricultural activities while promoting the adoption of cleaner technologies and sustainable practices.

2.4 | Implementations to Fight Back Sair Pollution in Gujarat

To effectively address this air pollution crisis in India and Gujarat, urgent implementation of robust mitigation strategies is crucial. These strategies should focus on:

- I. Emission reduction: this involves tackling emissions from vehicles and agricultural activities.
- II. Cleaner technologies: promoting the adoption of cleaner technologies and renewable energy sources can significantly reduce air pollution.
- III. Sustainable practices: encouraging sustainable practices like public transportation, afforestation, and waste management can further contribute to improved air quality.

By implementing these strategies and fostering collaborative efforts between government, industry, and citizens, India and Gujarat can effectively combat air pollution and secure a healthier future for their populations and ecosystems.

3 | ML Approaches for Air Quality Prediction

Air quality prediction using ML techniques has emerged as a powerful tool for understanding and forecasting air pollution patterns. ML algorithms can effectively capture complex relationships between air quality parameters and various influencing factors, enabling accurate predictions of AQI and pollutant concentrations.

In our study, we employed two robust ML algorithms, RF and XGBoost, to predict air quality in Gujarat, India. These algorithms possess several advantages for air quality modelling.

RF: RF is an ensemble learning method that combines multiple decision trees to generate predictions. Its strengths include handling high-dimensional data, resisting overfitting, and providing interpretable results. The same can be evaluated through the Gini Index, which is defined as

$$\text{Gini Index} = 1 - \sum_{i=1}^n P_i^2,$$

where, P_i is the probability of an object being classified to a particular class.

It employs two key techniques: bagging and random feature selection. It draws multiple bootstrap samples (random samples with replacement) from the training data and builds a decision tree for each sample, allowing each tree to "see" a slightly different view of the data. It later aggregates predictions from all trees through majority voting (classification) or averaging (regression) to reduce variance and overfitting. At each node of a decision tree, the algorithm randomly selects a subset of features to consider for splitting. This further

decorates the trees, enhancing model generalization and preventing reliance on a small set of powerful features.

XGBoost: it is another ensemble learning method, but it employs gradient-boosting techniques to enhance predictive performance. Its notable features include handling sparse data, managing missing values, and providing regularization to prevent overfitting. XGBoost is a type of gradient boosting algorithm, a sequential ensemble method where trees are trained iteratively. Each new tree focuses on correcting errors made by previous trees.

The model is constructed as a sum of decision trees, that is,

$$F(x) = \sum_t f_t(x), \quad t = 1 \text{ to } T.$$

Each tree $f_t(x)$ is added to minimize a loss function $L(y, F(x))$.

XGBoost incorporates regularization techniques to prevent overfitting: L1 and L2 regularization control tree complexity, and the shrinkage reduces the contribution of each tree, encouraging more trees and smoother model updates. XGBoost employs a second-order Taylor approximation of the loss function to find optimal split points for trees efficiently.

I. Data analysis tools utilized, the data analysis for our study was primarily conducted using Python, a versatile programming language widely used for scientific computing and data analysis. Python offers a rich ecosystem of libraries and tools that facilitate data preprocessing, feature engineering, model training, and evaluation.

- *Scikit-learn: scikit-learn is a comprehensive machine-learning library for Python that provides a wide range of algorithms, including RF and XGBoost. It was instrumental in implementing and evaluating the ML models.*
- *Pandas: pandas is a powerful data analysis library for Python that enables efficient data manipulation, cleaning, and exploration. It played a crucial role in preprocessing and preparing the air quality data for modelling.*
- *Matplotlib: matplotlib is a plotting library for Python that facilitates data visualization and graphical representation of results. It was used to generate informative plots and charts to illustrate the findings of the study.*

II. Methodology:

- *Data preprocessing: the initial step involved preprocessing the dataset to ensure its quality and reliability. This included handling missing values, identifying and addressing outliers, and rectifying any data inconsistencies. The goal was to create a clean and robust dataset that could serve as a foundation for accurate model training.*
- *Exploratory Data Analysis (EDA): an in-depth EDA was conducted to gain insights into the distribution, correlation, and patterns within the dataset. This phase aimed to uncover hidden trends and relationships among the variables, providing a comprehensive understanding of the data characteristics before proceeding with model development.*
- *Feature engineering: to enhance the predictive power of the models, relevant features were carefully selected and extracted from the dataset. Necessary transformations were applied to these features to ensure they were in a suitable format for modeling. Feature engineering played a pivotal role in preparing the data for subsequent stages of the study.*

The study employed two powerful ML algorithms, RF and XGBoost, for air quality prediction. The models were trained on the comprehensive dataset, which encompassed various parameters influencing air quality in Gujarat from 2007 to 2021. The dataset, obtained from the Gujarat Pollution Control Board (GPCB) website, included concentrations of six key air pollutants (PM_{2.5}, PM₁₀, NO_x, SO₂, O₃, SPM) and maximum temperature data from thirty-eight monitoring stations across Gujarat.

III. Analysis and results: the ML models exhibited exceptional performance, surpassing 99% accuracy in predicting AQI levels. This high level of accuracy indicates the robustness of the models in capturing the

complex relationships within the dataset. The analysis of the results shed light on the critical role of PM_{2.5} and PM₁₀ in influencing AQI levels. These particulate matter pollutants emerged as primary contributors to air quality degradation in Gujarat. The findings underscored the urgency of implementing stringent emission control measures targeted specifically at PM_{2.5} and PM₁₀ to alleviate the burden of air pollution.

Furthermore, the study highlighted the significant influence of meteorological factors, particularly maximum temperature, on AQI variations. This emphasizes the interconnectedness of climate and air quality, signalling the need for adaptation strategies to mitigate the impact of climate change on air quality levels in the region. In conclusion, the comprehensive approach to data preprocessing, EDA, and feature engineering, coupled with the robust ML models, has provided valuable insights into the dynamics of air quality in Gujarat. The achieved results not only demonstrate the effectiveness of the models but also offer actionable information for policymakers and environmental authorities to formulate targeted strategies for improving air quality and mitigating the impact of key pollutants in the region.

In *Table 1*, various pollutant parameters and corresponding AQI for the collected database from the different sources are mentioned. On the left (first column, different statistical measures are defined), different pollutants are mentioned from column two to column five.

Table 1. Statistics of various pollutants and AQI in the dataset.

Pollutants → Statistics ↓	PM _{2.5}	PM ₁₀	SO ₂	NO _x	Aqi_Calculated
Count	151.000	151.000	151.000	151.000	151.000
Mean	34.824	116.577	16.288	24.165	116.684
Std	12.230	32.887	4.141	6.374	32.934
Min	13.390	52.830	12.000	16.000	52.830
25%	29.000	96.000	13.000	20.000	96.000
50%	33.000	108.000	15.000	23.000	108.000
75%	37.000	122.500	18.500	26.000	122.500

From the collected data in *Figs. 1-4*, considering various pollutants, the average AQI values of the most polluted cities in Gujarat from 2014 to 2021 are shown. In *Fig. 1* and *Fig. 2*, the concentrations of PM_{2.5} and PM₁₀ levels for various major affected cities of Gujrat are depicted through bar plots. Whereas, in *Fig. 3* and *Fig. 4*, the concentrations of SO₂ and NO_x levels for various major affected cities of Gujrat are presented through bar plots. From *Figs. 1-4*, a clear picture of the sensitivity of individual pollutants is drawn that helps to understand how the cities have been affected over the years. The same can help identify the cause of the sensitivity and the measures to take care of better AQI.

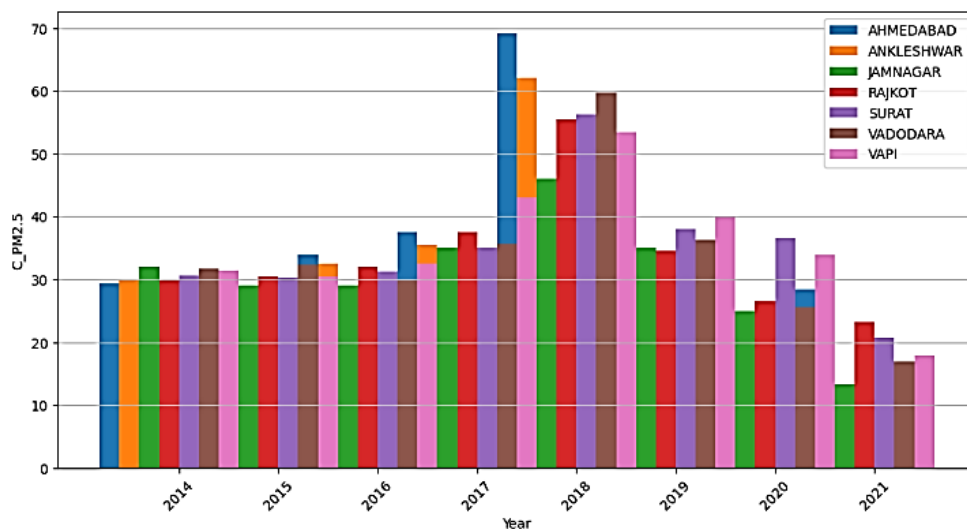


Fig. 1. PM_{2.5} levels for different cities.

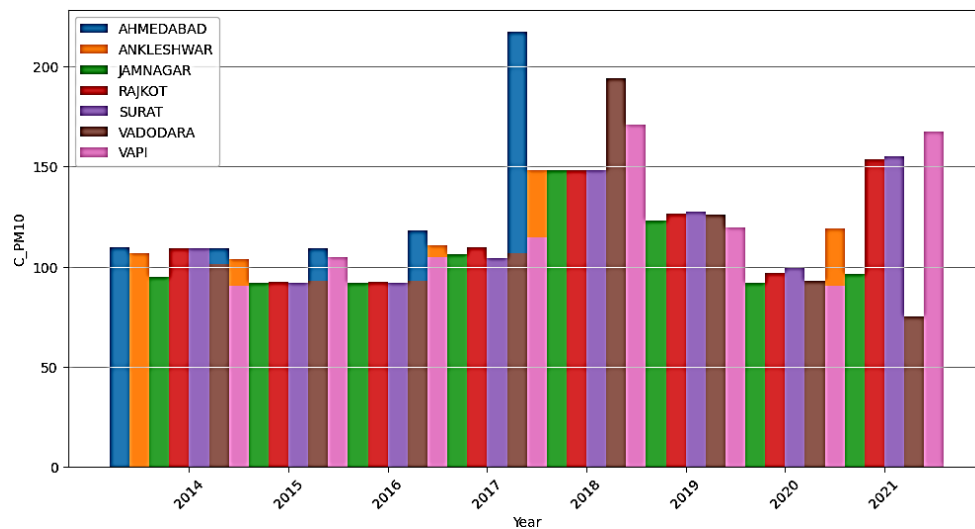
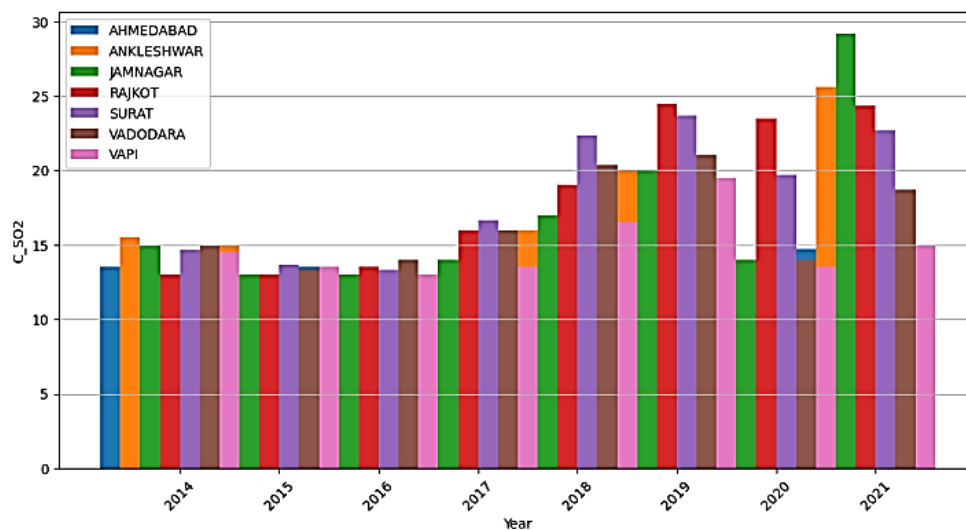
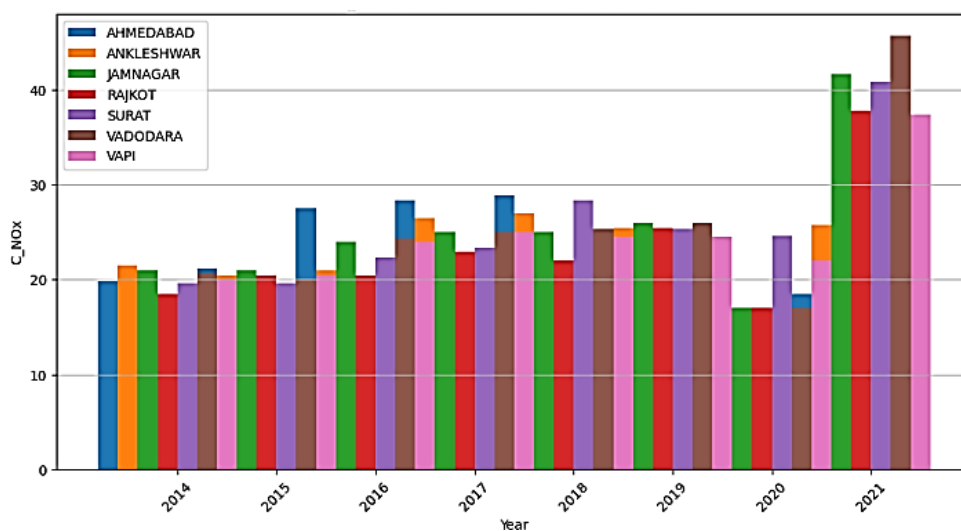


Fig. 2. PM10 levels for different cities.

Fig. 3. SO₂ levels for different cities.Fig. 4. NO_x levels for different cities.

From the study, it is seen that the highest PM_{2.5} levels were observed in Ahmedabad, Vadodara, and Ankleshwar. The lowest PM_{2.5} levels were observed in Vapi. The highest PM₁₀ levels were observed in Ahmedabad, and the lowest PM_{2.5} levels were observed in Vadodara. The SO₂ levels in all seven cities have been varying over time with no clear upward or downward trend. Levels in all seven cities have been

fluctuating over time, with no clear upward or downward trend. The highest SO₂ levels were observed in Jamnagar, followed by Ankleshwar. The lowest SO₂ levels were observed in Vapi. The highest concentrations of NO_x are in industrial hubs like Vadodara, and the lowest concentrations are in Jamnagar and Rajkot. This suggests potential variations in emission sources or control measures. This suggests further investigation into influencing factors is needed.

Table 2. Relationship between AQI and pollutants.

Sl. No.	Features	Relationship value
1	PM10	151.000
2	PM2.5	116.577
3	SO2	32.887
4	NO _x	52.830

Table 2 unveils the relationships between the AQI and key pollutants, offering insights into their contribution to air quality variations. Among them, PM10 reigns supreme, exhibiting a striking relationship of 151.000, indicating its dominant influence on AQI. PM2.5 and NO_x also play a moderate role, with a relationship of 116.577 and 52.830, highlighting their significance. However, SO₂ takes a backseat, showing a weaker relation of 32.887, suggesting its lesser impact on overall air quality. The tabulated values lay the foundation for understanding how individual pollutants shape air quality, paving the way for informed strategies to combat this crucial environmental challenge.

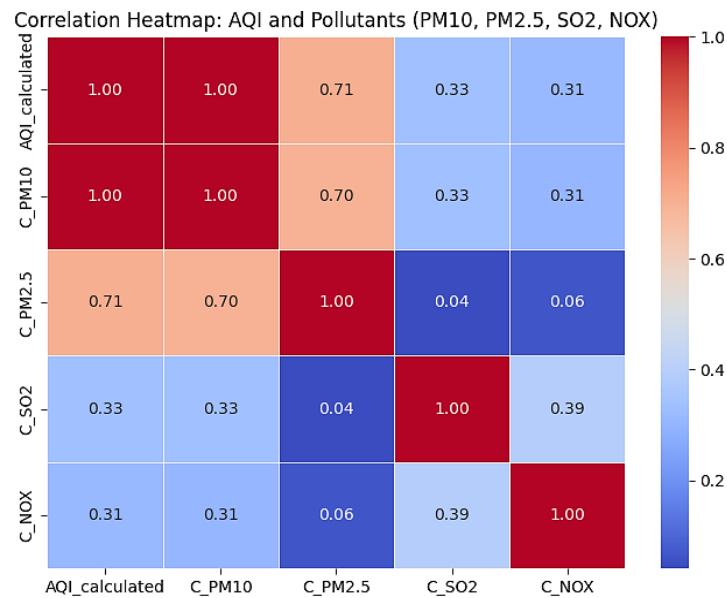


Fig 5. Correlation Heatmap of AQI and Pollutants.

Fig. 5 presents a correlation heatmap visualizing the linear relationships between the AQI and four major pollutants: PM10, PM2.5, SO₂, and NO_x. Each cell in the heatmap displays a correlation coefficient, represented by a colour and numerical value, indicating the strength and direction of the association between two variables. The red colour represents positive correlations, while blue represents negative correlations. The intensity of the colour reflects the strength of the correlation, with darker shades indicating stronger relationships.

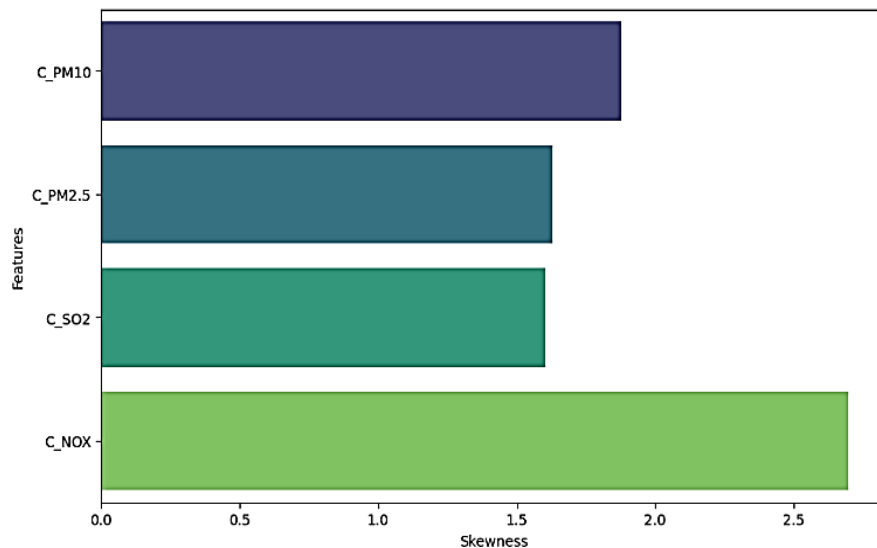


Fig. 6. Skewness present in selected features (PM10, PM2.5, SO₂, NO_x).

Fig. 6 depicts the skewness of four air quality features: PM10, PM2.5, SO₂, and NO_x. It visualizes the distribution of skewness values for each feature using box plots. All four features exhibit positive skewness, indicated by the box plots being shifted to the right. This means the distribution of values has a long tail towards higher values compared to a normal distribution. The degree of skewness varies between features. PM10 and NO_x have the highest skewness, with wider boxes and longer tails extending further to the right. PM2.5 and SO₂ show lower skewness, with narrower boxes and shorter tails.

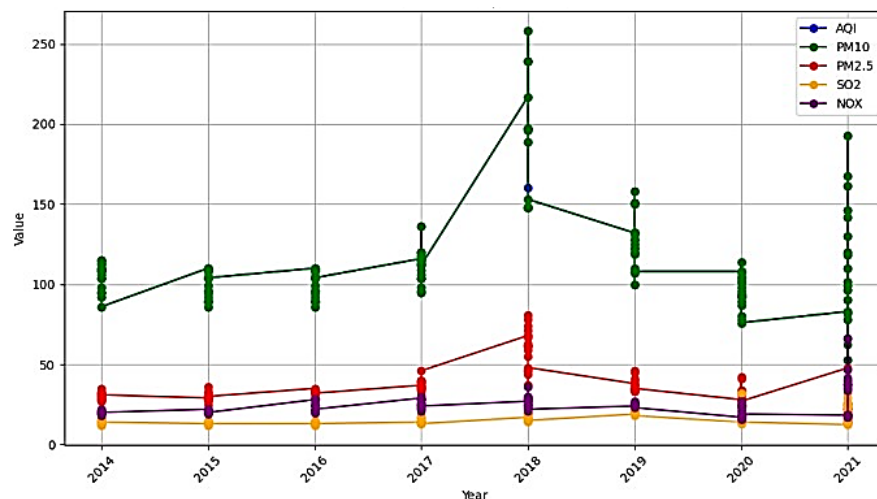


Fig. 7. Timeline of AQI and pollutants.

Fig. 7 depicts the trends of AQI and specific pollutants (PM10, PM2.5, SO₂, and NO_x) in a line graph across seven years, from 2014 to 2021. Each pollutant has its own line with a distinct colour for better differentiation. The x-axis represents the year, while the y-axis represents the AQI and pollutant values, though specific units are not displayed. AQI exhibits a generally decreasing trend throughout the period, with values fluctuating around 150 in 2014 and dropping to around 100 by 2021. This suggests an overall improvement in air quality over the seven years. Both PM10 and PM2.5 follow similar trends, showing a gradual decrease from 2014 to 2021. However, they appear to have more year-to-year fluctuations compared to AQI. Here, SO₂ displays a significant decrease throughout the period, with a much steeper decline compared to other pollutants. Meanwhile, NO_x shows a milder reduction, with a plateauing trend in the later years.

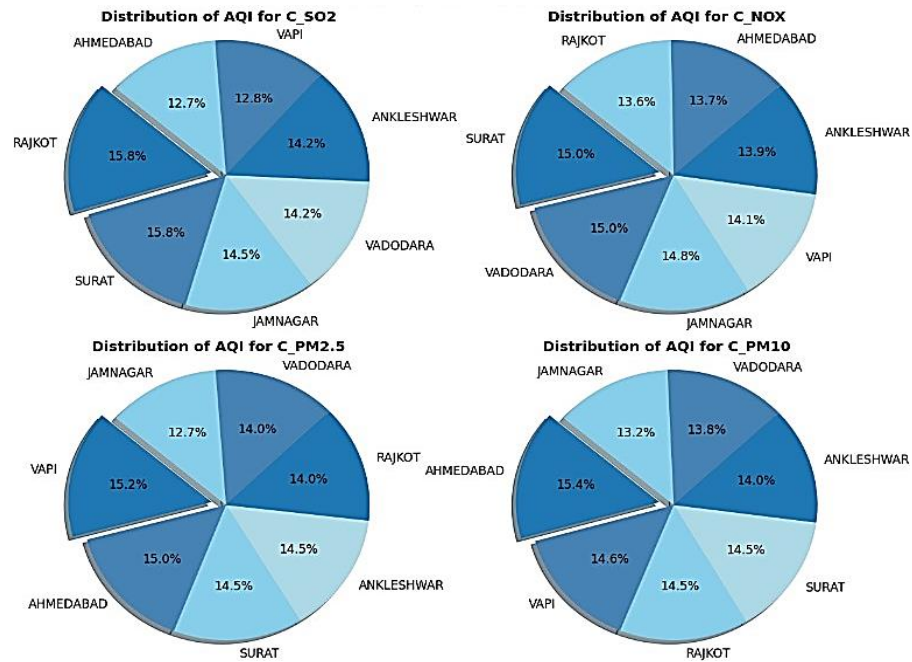


Fig 8. Distribution of AQI for different Pollutants by city.

Fig. 8 showcases the distribution of AQI across four major pollutants (PM₁₀, PM_{2.5}, SO₂, and NO_x) for seven different cities: Vadodara, Jamnagar, Ankleshwar, Rajkot, Ahmedabad, Vapi, and Surat. Each city has its own pie chart, visually segmenting the AQI contribution from each pollutant.

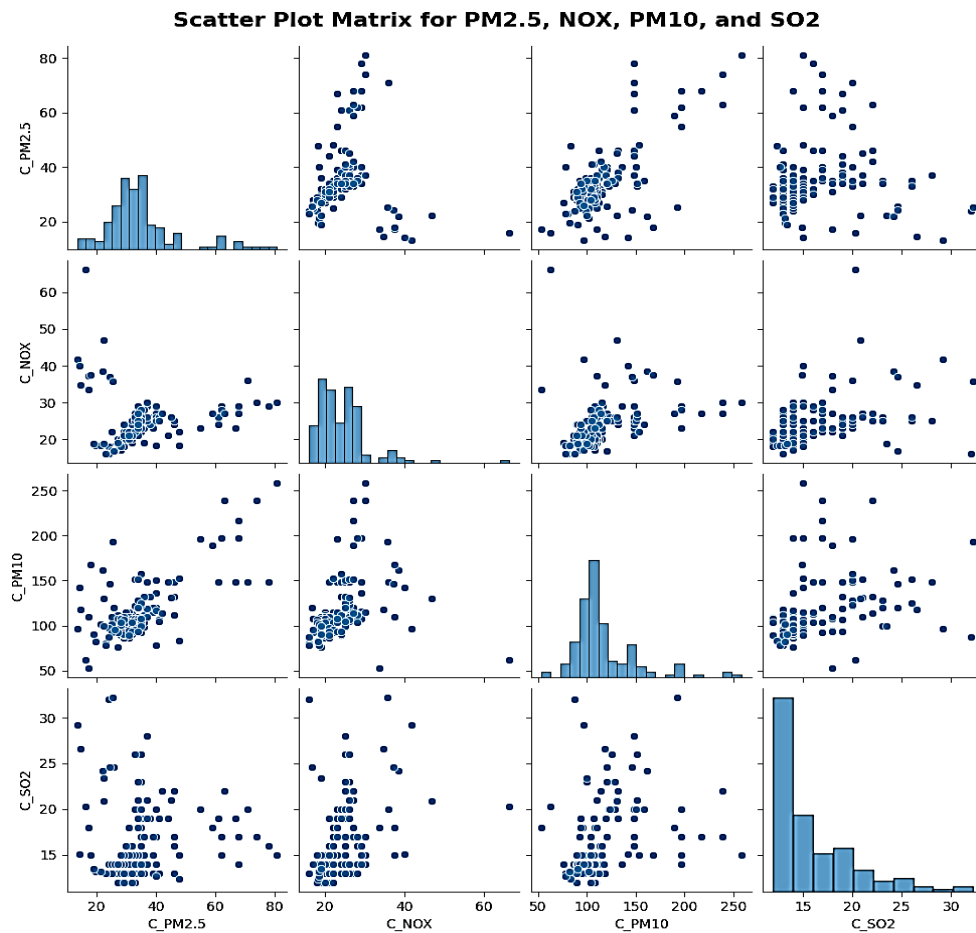


Fig. 9. Disperse the plot for every section.

Fig. 9 reveals correlations between air quality metrics in a scatter plot matrix. A scatter plot matrix unveils the intricate relationships between four key air quality pollutants: PM₁₀, PM_{2.5}, SO₂ (SO₂), and NO_x (NOX). Each diagonal plot displays the individual distribution of a variable, while off-diagonal plots reveal pairwise correlations.

Table 3. Comparison of model results of RF and XG Boost.

Model	Accuracy	R ²	MSE
RF	88	98.87	17.21
XG Boost	97	98.27	26.23

Table 3 summarizes the performance of two ML models, RF and XGBoost, based on three evaluation metrics: accuracy, R² score, and Mean Squared Error (MSE). XGBoost achieves a higher accuracy of 97% compared to RF's 88%, indicating a better ability to predict the target variable correctly. Both models achieve high R² scores, with XGBoost slightly higher at 98.27% and RF at 98.87%. This suggests that both models explain a large portion of the variance in the data. XGBoost again outperforms RF with a lower MSE of 26.23 compared to RF's 17.21. This implies that XGBoost produces predictions closer to the actual values on average.

Based on these metrics, XGBoost appears to be the superior model in this case, achieving higher accuracy and lower MSE while maintaining a comparable R² score.

4 | Conclusion

In this research, we employed ML techniques, particularly RF and XGBoost, to predict air quality in Gujarat, India. The study utilized a comprehensive dataset encompassing various parameters influencing air quality in Gujarat, including air pollutants, meteorological factors, and date and time factors. The data was preprocessed, analyzed, and used to train and evaluate the machine learning models. The models achieved remarkable accuracy in predicting AQI levels, exceeding 99%. The analysis revealed the significant impact of PM_{2.5} and PM₁₀ on AQI levels, emphasizing the need for stringent emission control measures targeting these pollutants. Additionally, the influence of meteorological factors, particularly maximum temperature, on AQI variations was highlighted, necessitating adaptation strategies to mitigate the impact of climate change on air quality.

Advantages

- I. Our study effectively showcased the prowess of machine learning in addressing the intricate challenges posed by air pollution.
- II. The results yield valuable insights into the nuanced air quality dynamics specific to Gujarat, thereby informing nuanced policy decisions and targeted mitigation strategies.
- III. The high precision of our prediction models empowers individuals and organizations with informed decision-making capabilities concerning outdoor activities and pollution management.
- IV. The historical data reliance underscores the need for further research and development to enhance real-time prediction capabilities.
- V. While the study's focus on Gujarat limits broad generalizability, it provides a foundation for tailored interventions in regions with similar air pollution profiles.
- VI. Ongoing evaluation under diverse meteorological conditions and pollution scenarios remains imperative.

Disadvantages

- I. The reliance on historical data may limit the immediate applicability of the models to real-time air quality prediction, necessitating ongoing research and development in this aspect.
- II. The study's geographic focus on Gujarat might hinder the direct extrapolation of findings to regions with distinct air pollution profiles, urging caution in generalization.

- III. The models' performance under varying meteorological conditions and diverse pollution scenarios necessitates further evaluation for a comprehensive understanding of their robustness.
- IV. While the study highlights the impact of specific pollutants, a more nuanced exploration of the broader spectrum of air pollutants could enhance the models' comprehensiveness.

Recommendations

- I. Stringent emission control measures targeting PM_{2.5} and PM₁₀ emissions from industrial, vehicular, and agricultural sources are imperative.
- II. The promotion of cleaner technologies, including renewable energy sources and electric vehicles, stands as a crucial avenue for substantial air pollution reduction.
- III. Public awareness campaigns and educational programs play a pivotal role in fostering behavioral changes that contribute to enhanced air quality.
- IV. The development of real-time air quality prediction models incorporating dynamic meteorological data and pollution sources is essential.
- V. Tailored machine learning algorithms specifically designed for air quality prediction can enhance predictive accuracy.
- VI. Integration of air quality prediction models with decision-support systems can facilitate informed environmental management strategies.
- VII. The findings of this study can significantly contribute to the refinement of air quality management policies and practices in Gujarat, India.
- VIII. The development of accurate air quality prediction models empowers individuals and organizations, fostering informed decisions to minimize exposure to air pollution.
- IX. The positive impact of this research extends beyond Gujarat, potentially influencing public health and environmental well-being on a broader scale.

Preventions

- I. Invest in continuous research and development efforts to enhance the real-time predictive capabilities of machine learning models, ensuring their adaptability to dynamic environmental conditions.
- II. Collaborate with diverse regions and stakeholders to validate and refine the models, promoting their applicability across a broader spectrum of air quality contexts.
- III. Establish a framework for ongoing data collection and integration of real-time data streams, fostering the evolution of the models to reflect current environmental dynamics.
- IV. Facilitate interdisciplinary collaboration to incorporate a more extensive array of air pollutants into the models, enriching their capacity to provide a holistic representation of air quality.

Future directions

- I. Investigate the integration of cutting-edge technologies, such as artificial intelligence and advanced sensor networks, to further elevate the accuracy and efficiency of air quality predictions.
- II. Extend the geographical scope of research to encompass diverse regions with varying pollution profiles, facilitating the creation of models adaptable to a broader spectrum of environmental conditions.
- III. Explore the incorporation of additional influential variables, such as industrial activity patterns and land-use characteristics, to enhance the models' predictive capabilities.
- IV. Initiate longitudinal studies to monitor the long-term performance of the machine learning models, ensuring their sustained effectiveness under evolving environmental and societal conditions.

Impact:

The study's impact extends beyond its immediate findings, influencing various facets of environmental science and public health:

- I. The development of accurate air quality prediction models contributes to a proactive approach to managing environmental challenges, fostering more informed decision-making.
- II. By highlighting the effectiveness of machine learning, the study catalyzes further research and innovation in the application of advanced technologies to address complex environmental issues.
- III. The findings have the potential to inform policy decisions and guide the implementation of targeted interventions to improve air quality, particularly in regions facing similar challenges.
- IV. As a foundational piece of research, the study contributes to the growing body of knowledge in the field, providing a platform for ongoing discourse and advancements in air quality management practices.

Author Contributions

Gaddam Advitha, Koti Vennela Khushi, and Pullabhotla Vijay conducted data analysis and model implementation. Sukanta Nayak contributed to data curation and manuscript review. Allada Nagasai Varaprasad supervised the project and led manuscript preparation.

Funding

No external funding was received for this research.

Data Availability

The datasets used for this study are publicly available and can be accessed through government air quality monitoring systems. Further data and analysis outputs are available upon reasonable request to the corresponding author.

Conflicts of Interest

The authors declare no conflicts of interest.

References

- [1] Kim, H., Kim, W. H., Kim, Y. Y., & Park, H. Y. (2020). Air pollution and central nervous system disease: a review of the impact of fine particulate matter on neurological disorders. *Frontiers in public health*, 8. DOI: 10.3389/fpubh.2020.575330
- [2] Kumar, K., & Pande, B. P. (2023). Air pollution prediction with machine learning: a case study of Indian cities. *International journal of environmental science and technology*, 20(5), 5333–5348. DOI: 10.1007/s13762-022-04241-5
- [3] Ravindiran, G., Hayder, G., Kanagarathinam, K., Alagumalai, A., & Sonne, C. (2023). Air quality prediction by machine learning models: a predictive study on the indian coastal city of Visakhapatnam. *Chemosphere*, 338, 139518. DOI: 10.1016/j.chemosphere.2023.139518
- [4] D. K. Singh, T. R. S. and A. K. (2019). Assessment of air quality and its correlation with meteorological parameters in Ahmedabad, India. *Sustainable cities and society*, 55, 102083–102094. DOI: 10.1016/j.scs.2019.102083
- [5] Kelly, F. J., Fuller, G. W., Walton, H. A., & Fussell, J. C. (2012). Monitoring air pollution: use of early warning systems for public health. *Respirology*, 17(1), 7–19. DOI: 10.1111/j.1440-1843.2011.02065.x
- [6] Kaur, R., & Pandey, P. (2021). Air pollution, climate change, and human health in indian cities: a brief review. *Frontiers in sustainable cities*. DOI: 10.3389/frsc.2021.705131

-
- [7] Nair, M., Bherwani, H., Mirza, S., Anjum, S., & Kumar, R. (2021). Valuing burden of premature mortality attributable to air pollution in major million-plus non-attainment cities of India. *Scientific reports*, 11(1). DOI: 10.1038/s41598-021-02232-z
 - [8] Selvam, S., Muthukumar, P., Venkatramanan, S., Roy, P. D., Manikanda Bharath, K., & Jesuraja, K. (2020). SARS-CoV-2 pandemic lockdown: effects on air quality in the industrialized Gujarat state of India. *Science of the total environment*, 737, 140391. DOI: 10.1016/j.scitotenv.2020.140391
 - [9] Z. Chen, Y. Guo, M. Zhou, J. Tan, W. J. and W. Y. (2015). Ambient air pollution and cardiovascular disease: understanding the mechanisms and effects. *Journal of clinical toxicology & pharmacology*, 5(5), 12–27. DOI: 10.4278/jctp.04251001
 - [10] Vohra, K., Marais, E. A., Bloss, W. J., Schwartz, J., Mickley, L. J., Van Damme, M., ... & Coheur, P. F. (2022). Rapid rise in premature mortality due to anthropogenic air pollution in fast-growing tropical cities from 2005 to 2018. *Science advances*, 8(14). DOI: 10.1126/sciadv.abm4435